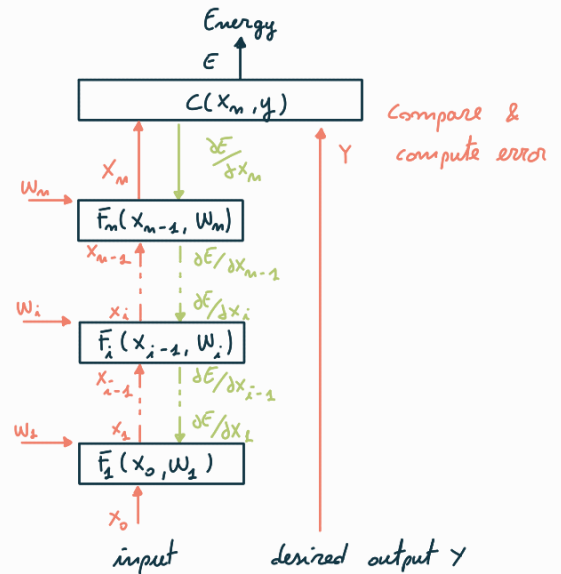# Gradient computation

Information flows forward through the network when computing network output y from input x

$\longrightarrow$

The Backpropagation algorithm is used to propagate gradient computation from the cost function through the whole network

$\longleftarrow$

Energy
$E$

$C(x_m, y)$ — Compare & compute error

$w_m$   $x_m$   $\delta E/\delta x_m$   $Y$

$F_m(x_{m-1}, w_m)$

$x_{m-1}$   $\delta E/\delta x_{m-1}$

$w_i$   $x_i$   $\delta E/\delta x_i$

$F_i(x_{i-1}, w_i)$

$x_{i-1}$   $\delta E/\delta x_{i-1}$

$w_1$   $x_1$   $\delta E/\delta x_1$

$F_1(x_o, w_1)$

$x_o$

input          desired output $Y$

Goal: compute the gradient of the cost function with respect to the parameters:

$$\nabla_\theta J(\theta)$$

Backpropagation is not a training algorithm, it is just an algorithm to estimate the gradient. Training will be done by an optimizer (e.g. Adam) by means of the stochastic gradient descent.

## Chain rule

Let $y = g(x)$   and   $z = f(g(x)) = f(y)$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

For vector functions $g: R^m \longrightarrow R^n$ and $f \cdot R^n \longrightarrow R$ we have:
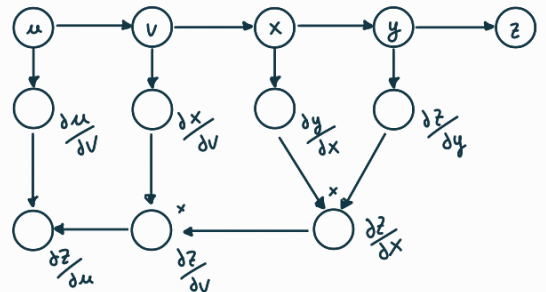
$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

In vector notation:

$$\nabla_x z = \left(\frac{\partial y}{\partial x}\right)^T \nabla_y z$$

with $\frac{\partial y}{\partial x}$ the $n \times m$ Jacobian matrix of $g$

## Why is it so easy to compute the gradient?

$u \to v \to x \to y \to z$

$\frac{\partial u}{\partial v}$   $\frac{\partial x}{\partial v}$   $\frac{\partial y}{\partial x}$   $\frac{\partial z}{\partial y}$

$\frac{\partial z}{\partial x}$

$\frac{\partial z}{\partial u}$   $\frac{\partial z}{\partial v}$

Estimating gradient $\longrightarrow$ Backpropagation

Training $\longrightarrow$ training algorithms $\longrightarrow$ Stochastic Gradient Descent (SGD)
$\longrightarrow$ SGD with momentum
$\longrightarrow$ Algorithms with adaptive learning rates

# Backpropagation

Backpropagation consists of two steps : ① forward step and ② backward step

## ① Forward step

**Require:** Network depth $l$
**Require:** $W^{(i)}, i \in \{1, \ldots, l\}$ weight matrices
**Require:** $\mathbf{b}^{(i)}, i \in \{1, \ldots, l\}$ bias parameters
**Require:** $\mathbf{x}$ input value
**Require:** $\mathbf{t}$ target value

$h^{(0)} = \mathbf{x}$
**for** $k = 1, \ldots, l$ **do**
$\quad \boldsymbol{\alpha}^{(k)} = \mathbf{b}^{(k)} + W^{(k)} \mathbf{h}^{(k-1)}$
$\quad \mathbf{h}^{(k)} = f(\boldsymbol{\alpha}^{(k)})$
**end for**
$\mathbf{y} = \mathbf{h}^{(l)}$
$J = L(\mathbf{t}, \mathbf{y})$

The second step computes the derivative of the error with respect to all the parameters of the network, going backwards from the output layer to the input layer

## ② Backward step

*g is the portion of the gradient that has been computed so far*

① $\mathbf{g} \leftarrow \nabla_{\mathbf{y}} J = \nabla_{\mathbf{y}} L(\mathbf{t}, \mathbf{y})$
**for** $k = l, l-1, \ldots, 1$ **do**
$\quad$ Propagate gradients to the pre-nonlinearity activations:
② $\quad \mathbf{g} \leftarrow \nabla_{\boldsymbol{\alpha}^{(k)}} J = \mathbf{g} \odot f'(\boldsymbol{\alpha}^{(k)})$ {$\odot$ denotes elementwise product}
③ $\quad \nabla_{\mathbf{b}^{(k)}} J = \mathbf{g}$
④ $\quad \nabla_{W^{(k)}} J = \mathbf{g}(\mathbf{h}^{(k-1)})^T$
$\quad$ Propagate gradients to the next lower-level hidden layer:
⑤ $\quad \mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = (W^{(k)})^T \mathbf{g}$
**end for**

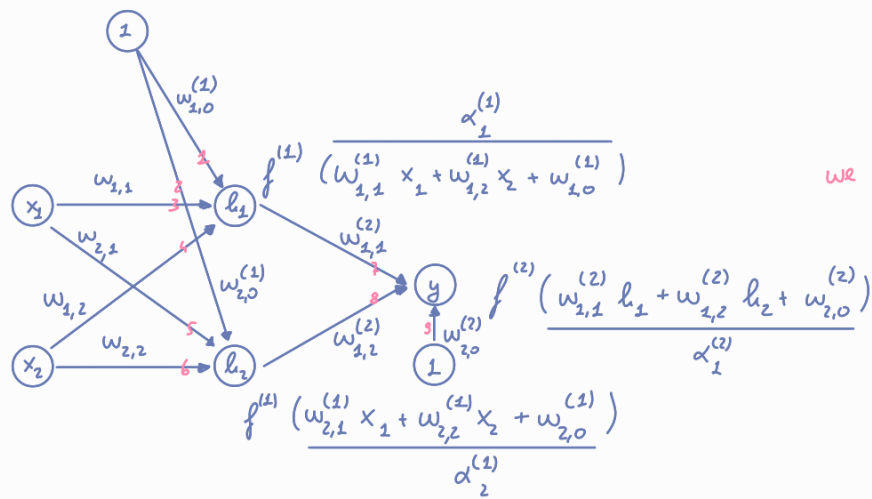*at each layer we need to multiply the contribution of $f(\alpha)$ : $f'(\alpha) \cdot$ derivative $(\alpha)$*

$\alpha_i = w_i^T h_{i-1} + b_i$
$\frac{\partial \alpha_i}{\partial w_i} = h_{i-1}$  $\frac{\partial \alpha_i}{\partial b_i} = 1$

This is not the training, it is just a part of it

Example



$$f^{(1)} \frac{(w_{1,1}^{(1)} x_1 + w_{1,2}^{(1)} x_2 + w_{1,0}^{(1)})}{\alpha_1^{(1)}}$$

$$f^{(2)} \frac{(w_{1,1}^{(2)} h_1 + w_{1,2}^{(2)} h_2 + w_{2,0}^{(2)})}{\alpha_1^{(2)}}$$

$$f^{(1)} \frac{(w_{2,1}^{(1)} x_1 + w_{2,2}^{(1)} x_2 + w_{2,0}^{(1)})}{\alpha_2^{(1)}}$$

we want to compute 9 derivatives

$$f^{(1)}(z) = ReLU$$
$$f^{(2)}(z) = z$$
$$Loss = MSE$$

Forward step
$\begin{cases} \text{Given } x_1, x_2, w_{i,j}^{(k)}, t \\ \text{Compute } \alpha_1^{(1)}, \alpha_2^{(1)}, \alpha^{(2)}, h_1^{(1)}, h_2^{(1)}, h_1^{(2)}, y, J = L(t,y) \end{cases}$

**layer 1**

$$\alpha_1^{(1)} = w_{1,0}^{(1)} + w_{1,1}^{(1)} x_1 + w_{1,2}^{(1)} x_2 \qquad h_1^{(1)} = f(\alpha_1^{(1)})$$

$$\alpha_2^{(1)} = w_{2,0}^{(1)} + w_{2,1}^{(1)} x_1 + w_{2,2}^{(1)} x_2 \qquad h_2^{(1)} = f(\alpha_2^{(1)})$$

$\begin{cases} \alpha_i^{(1)} = w_{i,0}^{(1)} + w_{i,1}^{(1)} x_1 + w_{i,2}^{(1)} x_2 \\ h_i^{(1)} = f^{(1)}(\alpha_i^{(1)}) \end{cases}$

**layer 2**

$$\alpha_1^{(2)} = w_{1,0}^{(2)} + w_{1,1}^{(2)} h_1 + w_{1,2}^{(2)} h_2 \qquad h_1^{(2)} = f(\alpha_1^{(2)}) \qquad y = h^{(2)}$$

Let's use the compact notation

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \qquad W^{(1)} = \begin{bmatrix} w_{1,1}^{(1)} & w_{1,2}^{(1)} \\ w_{2,1}^{(1)} & w_{2,2}^{(1)} \end{bmatrix} \qquad b^{(1)} = \begin{bmatrix} w_{1,0}^{(1)} \\ w_{2,0}^{(1)} \end{bmatrix} \qquad W^{(2)} = \begin{bmatrix} w_{1,1}^{(2)} & w_{1,2}^{(2)} \end{bmatrix} \qquad b^{(2)} = \begin{bmatrix} w_{1,0}^{(2)} \end{bmatrix}$$

$$h^{(0)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x$$

**layer 1**

$$\alpha^{(1)} = \begin{bmatrix} \alpha_1^{(1)} \\ \alpha_2^{(1)} \end{bmatrix} = W^{(1)} h^{(0)} + b^{(1)} \qquad h^{(1)} = \begin{bmatrix} h_1^{(1)} \\ h_2^{(1)} \end{bmatrix} = f^{(1)}(\alpha^{(1)}) = \begin{bmatrix} f^{(1)}(\alpha_1^{(1)}) \\ f^{(1)}(\alpha_2^{(1)}) \end{bmatrix} = \begin{bmatrix} ReLU(\alpha_1^{(1)}) \\ ReLU(\alpha_2^{(1)}) \end{bmatrix}$$

**layer 2**

$$\alpha^{(2)} = \begin{bmatrix} \alpha_1^{(2)} \end{bmatrix} = W^{(2)} h^{(1)} + b^{(2)} \qquad h^{(2)} = \begin{bmatrix} h_1^{(2)} \end{bmatrix} = f^{(2)}(\alpha^{(2)}) = \begin{bmatrix} \alpha_1^{(2)} \end{bmatrix}$$

$\begin{cases} \alpha^i = W^i h^{i-1} + b^i \\ h^i = f^i(\alpha^i) \end{cases}$

Compute loss function (MSE)  $L(t,y) = \frac{1}{2}(t-y)^2$

① $\quad g \longleftarrow \nabla_{h^{(2)}} J = \nabla_y J = \nabla_y \frac{1}{2}(t-y)^2 = \frac{1}{2} 2(t-y)^1(-1) = -(t-y) = y-t$

② $\quad g \longleftarrow \nabla_{\alpha^{(2)}} J = g \odot f^{(2)'}(\alpha^{(2)}) = g \odot \boxed{\dfrac{\partial \alpha^{(2)} - t}{\partial \alpha^{(2)}}} = g \odot 1 = g \qquad$ <span style="color:pink">linear activation function</span>

we computed the gradient of $J$ with respect to $\alpha^{(2)}$; now we want to know the influence of the weights and biases on $\partial J / \partial \alpha^{(2)}$; this will be useful later (SGD) to update them in a training algorithm.

**Layer 2**

③ $\quad \nabla_{b^{(2)}} J \longleftarrow \left[ \dfrac{\partial J}{\partial w^{(2)}_{1,0}} \right] = g$

④ $\quad \nabla_{W^{(2)}} J \longleftarrow \left[ \dfrac{\partial J}{\partial w^{(2)}_{1,1}} \quad \dfrac{\partial J}{\partial w^{(2)}_{1,2}} \right] = g \cdot (h^{(1)})^T$

$\qquad W^{(2)} = \left[ w^{(2)}_{1,1} \quad w^{(2)}_{1,2} \right]$

Propagate gradients to the next lower-level hidden layer

⑤ $\quad g \longleftarrow \nabla_{h^{(k-1)}} J = \nabla_{h^{(1)}} J = (W^{(k)})^T g = (W^{(2)})^T g = \left[ w^{(2)}_{1,1} \quad w^{(2)}_{1,2} \right]^T g$

**Layer 1**

② $\quad g \longleftarrow \nabla_{\alpha^{(1)}} J = g \odot f^{(1)'}(\alpha^{(1)}) = g \odot \left[ \begin{array}{c} \partial ReLU(\alpha^{(1)}_1)/\partial \alpha^{(1)}_1 \\ \partial ReLU(\alpha^{(1)}_2)/\partial \alpha^{(1)}_2 \end{array} \right] = g \odot \left[ \begin{array}{c} step(\alpha^{(1)}_1) \\ step(\alpha^{(1)}_2) \end{array} \right]$

③ $\quad \nabla_{b^{(1)}} J \longleftarrow \left[ \begin{array}{c} \partial J/\partial w^{(1)}_{1,0} \\ \partial J/\partial w^{(1)}_{2,0} \end{array} \right] = g$

④ $\quad \nabla_{W^{(1)}} J \longleftarrow \left[ \begin{array}{cc} \partial J/\partial w^{(1)}_{1,1} & \partial J/\partial w^{(1)}_{1,2} \\ \partial J/\partial w^{(1)}_{2,1} & \partial J/\partial w^{(1)}_{2,2} \end{array} \right] = g\,(h^{(1)})^T$